

Breaking the Sample Size Barrier in Reinforcement Learning

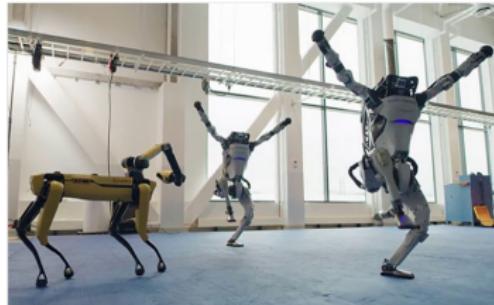


Yuting Wei

Carnegie Mellon University

Feb 2021

Reinforcement learning (RL)



In RL, an agent learns by interacting with an (unknown) environment.

Published: 25 February 2015

Human-level control through deep reinforcement learning

Volodymyr Mnih, Koray Kavukcuoglu✉, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg & Demis Hassabis✉

Nature **518**, 529–533(2015) | [Cite this article](#)

124k Accesses | **5073** Citations | **1539** Altmetric | [Metrics](#)

Published: 25 February 2015

Human-level control through reinforcement learning

Volodymyr Mnih, Koray Kavukcuoglu✉, David Bellemare, Alex Graves, Martin Riedmiller, And Charles Beattie, Amir Sadik, Ioannis Antonoglou, Shane Legg & Demis Hassabis✉

Nature 518, 529–533(2015) | Cite this article

124k Accesses | 5073 Citations | 1539 Alt

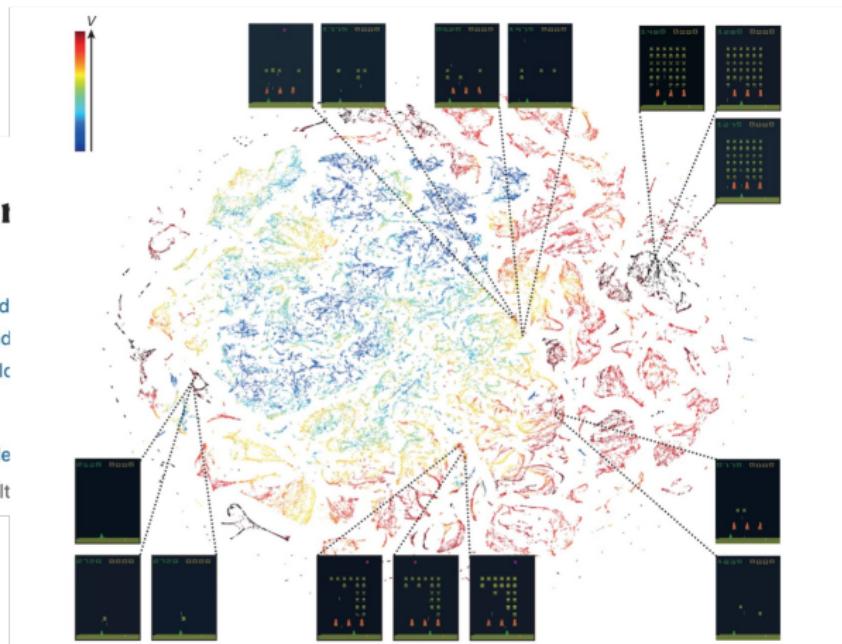
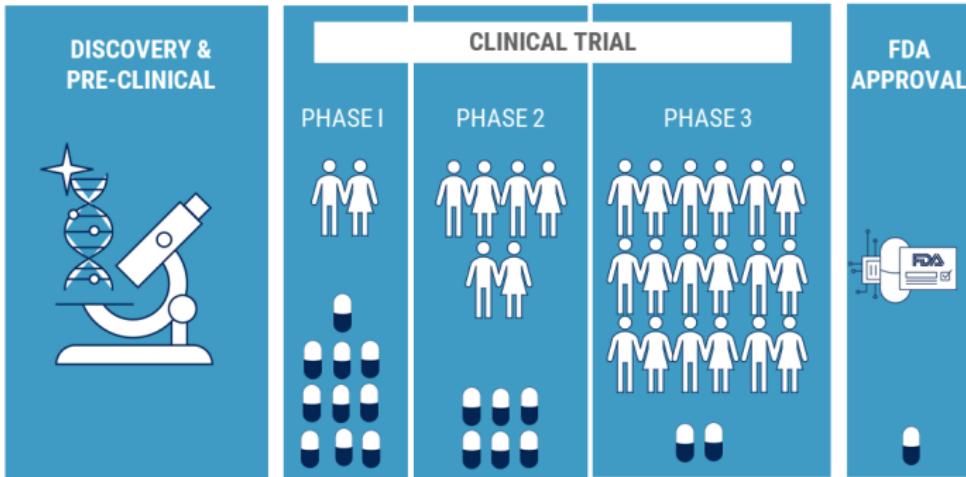


Figure 4: Two-dimensional t-SNE embedding of the representations in the last hidden layer assigned by DQN to game states experienced while playing Space Invaders.

From: Human-level control through deep reinforcement learning

Challenges in RL

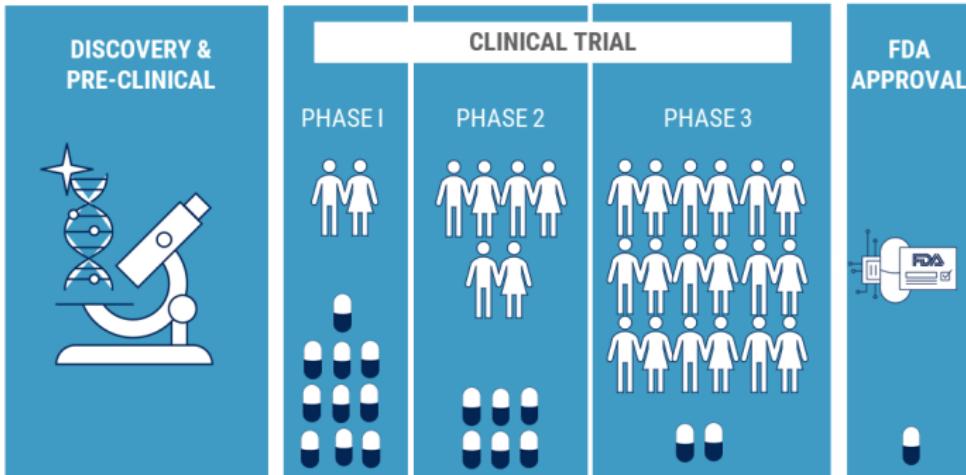


Source: cbinsights.com

CBINSIGHTS

- prohibitively large state & action space
- collecting data samples can be expensive or time-consuming

Challenges in RL



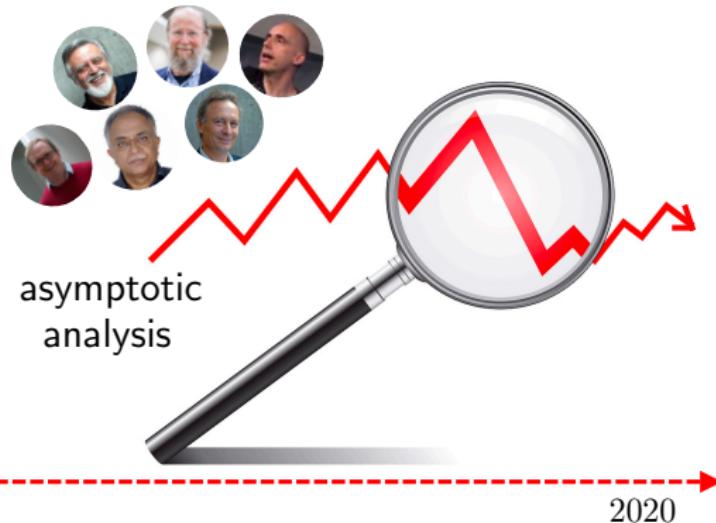
Source: cbinsights.com

CBINSIGHTS

- prohibitively large state & action space
- collecting data samples can be expensive or time-consuming

Challenge: design & understand sample efficient RL algorithms

Statistical foundation of RL



The Contributions of Herbert Robbins to Mathematical Statistics

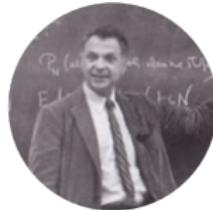
Tze Leung Lai and David Siegmund

2. STOCHASTIC APPROXIMATION AND ADAPTIVE DESIGN

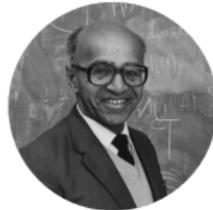
In 1951, Robbins and his student, Sutton Monroe, founded the subject of stochastic approximation with the publication of their celebrated paper [26]. Consider the problem of finding the root θ (assumed unique) of an equation $g(x) = 0$. In the classical

4. SEQUENTIAL EXPERIMENTATION AND OPTIMAL STOPPING

The well known “multiarmed bandit problem” in the statistics and engineering literature, which is prototypical of a wide variety of adaptive control and design problems, was first formulated and studied by Robbins [28]. Let A, B denote two statistical populations with finite means μ_A, μ_B . How should we draw a



Herbert Robbins



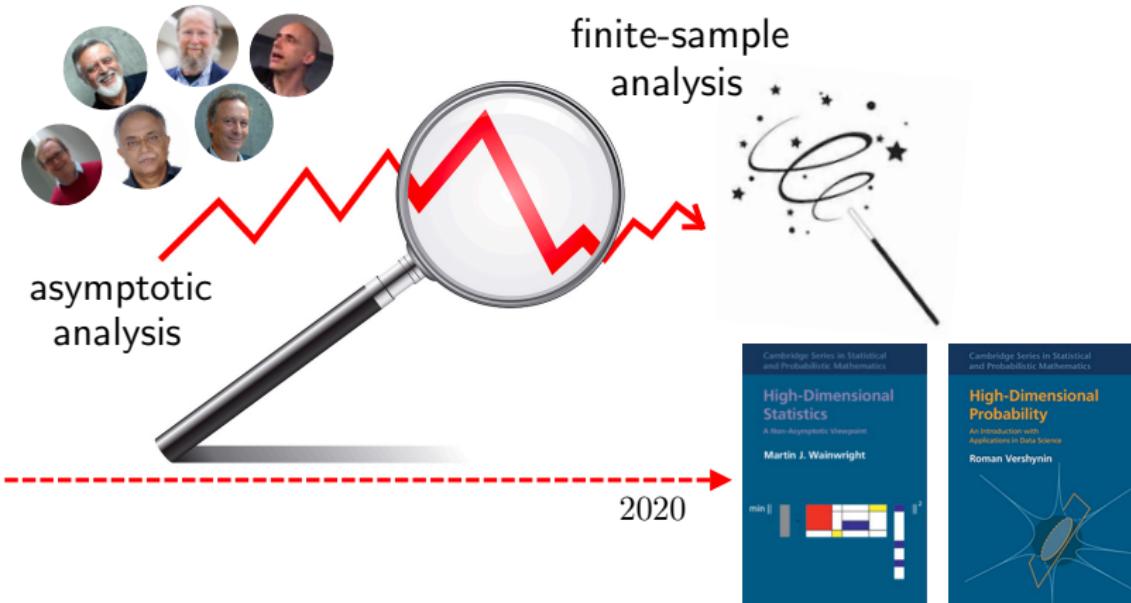
David Blackwell

David Blackwell, 1919–2010: An explorer in mathematics and statistics

Peter J. Bickel^{a,*†}

Blackwell channel. He also began to work in dynamic programming, which is now called reinforcement learning. In a series of papers, Blackwell gave a rigorous foundation to the theory of dynamic programming, introducing what have become known as Blackwell optimal policies.

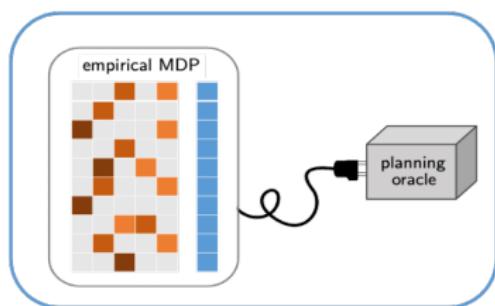
Statistical foundation of RL



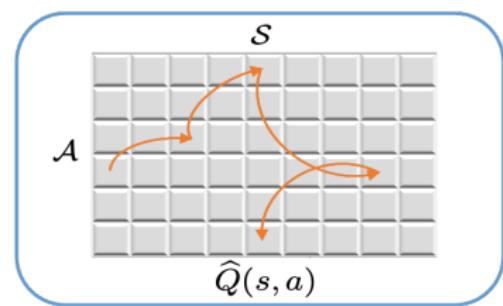
Understanding sample efficiency of RL requires a modern suite of non-asymptotic statistical tools.

Outline

- Background
- Vignette 1: model-based RL (“plug-in” approach)
- Vignette 2: model-free RL (Q-learning on Markovian samples)



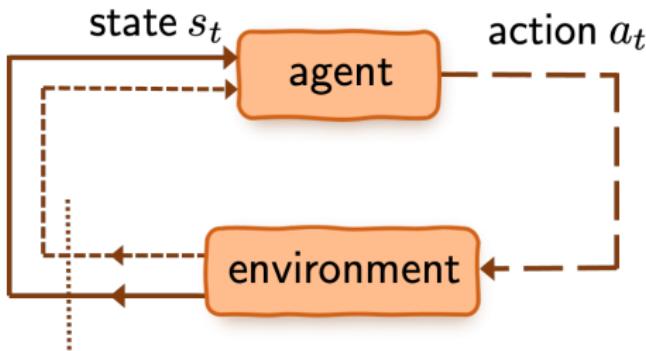
model based RL



model free RL

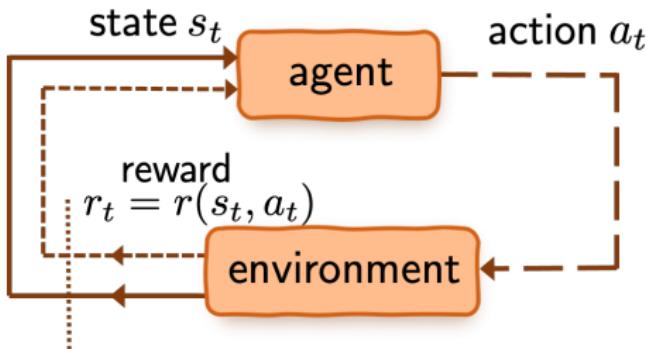
Background: Markov decision processes

Markov decision process (MDP)



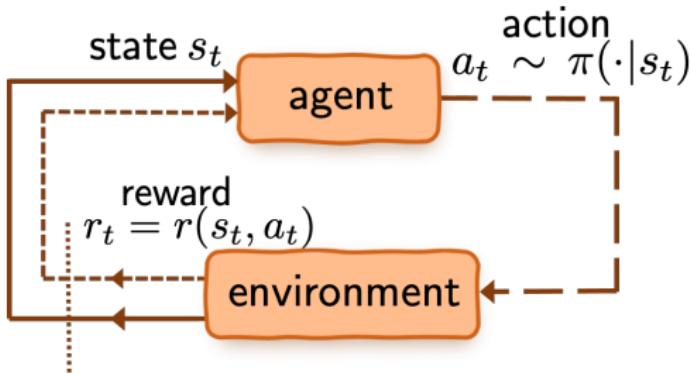
- \mathcal{S} : state space
- \mathcal{A} : action space

Markov decision process (MDP)



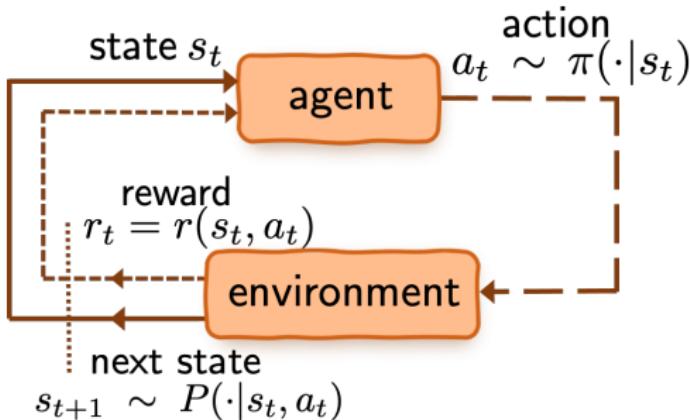
- \mathcal{S} : state space
- \mathcal{A} : action space
- $r(s, a) \in [0, 1]$: immediate reward

Markov decision process (MDP)



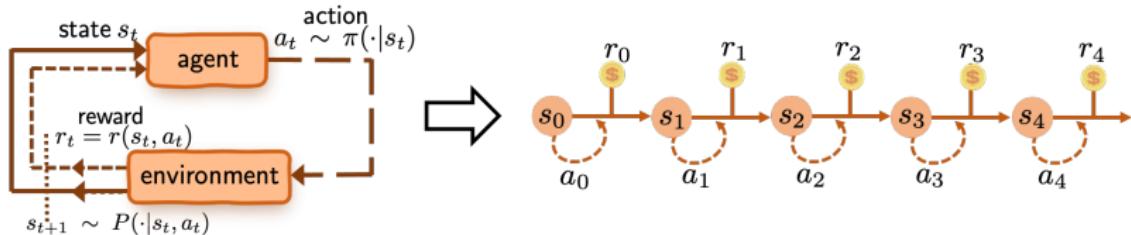
- \mathcal{S} : state space
- \mathcal{A} : action space
- $r(s, a) \in [0, 1]$: immediate reward
- $\pi(\cdot | s)$: policy (or action selection rule)

Markov decision process (MDP)



- \mathcal{S} : state space
- \mathcal{A} : action space
- $r(s, a) \in [0, 1]$: immediate reward
- $\pi(\cdot | s)$: policy (or action selection rule)
- $P(\cdot | s, a)$: **unknown** transition probabilities

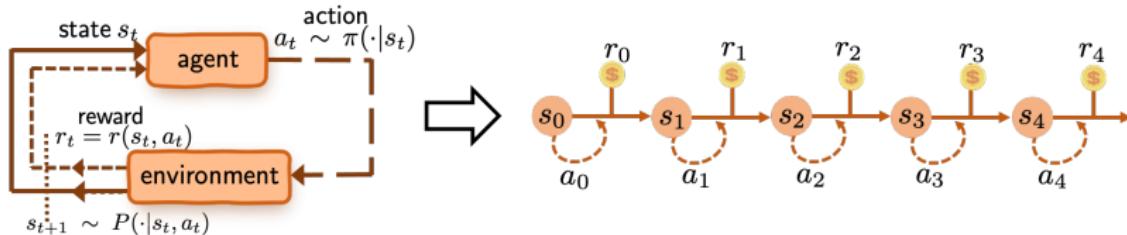
Value function



Value of policy π : cumulative **discounted** reward

$$\forall s \in \mathcal{S} : \quad V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$

Value function

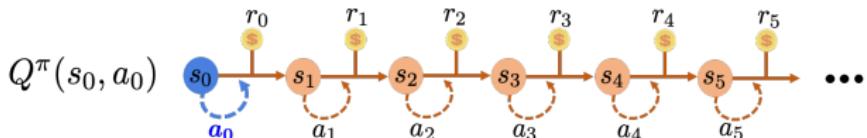


Value of policy π : cumulative **discounted** reward

$$\forall s \in \mathcal{S} : V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$

- $\gamma \in [0, 1)$: discount factor
 - ▶ take $\gamma \rightarrow 1$ to approximate **long-horizon** MDPs
 - ▶ **effective horizon**: $\frac{1}{1-\gamma}$

Q-function (action-value function)

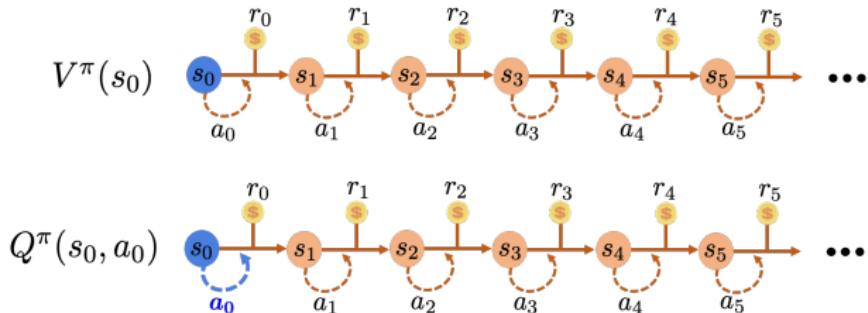


Q-function of policy π :

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A} : Q^\pi(s, a) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right]$$

- $(\cancel{a_0}, s_1, a_1, s_2, a_2, \dots)$: induced by policy π

Q-function (action-value function)



Q-function of policy π :

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}: \quad Q^\pi(s, a) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, \textcolor{red}{a_0 = a} \right]$$

- $(\textcolor{red}{a_0}, s_1, a_1, s_2, a_2, \dots)$: induced by policy π

Optimal policy and optimal value



- **optimal policy** π^* : maximizing value function $\max_{\pi} V^{\pi}(s)$

Optimal policy and optimal value



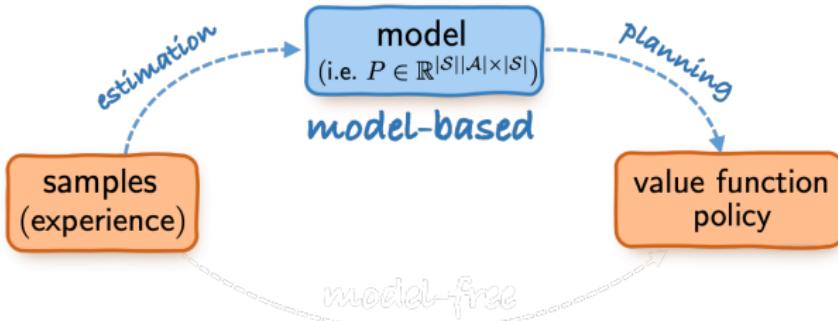
- **optimal policy** π^* : maximizing value function $\max_{\pi} V^{\pi}(s)$
- optimal value / Q function: $V^* := V^{\pi^*}$, $Q^* := Q^{\pi^*}$

Optimal policy and optimal value



- **optimal policy** π^* : maximizing value function $\max_{\pi} V^{\pi}(s)$
- optimal value / Q function: $V^* := V^{\pi^*}$, $Q^* := Q^{\pi^*}$
- How to find this π^* ?

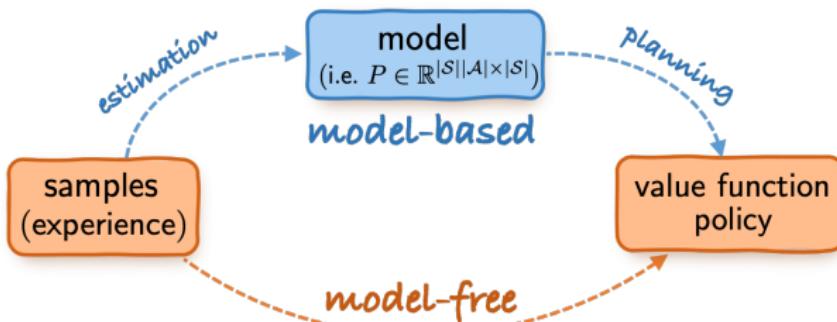
Model-based vs. model-free RL



Model-based approach (“plug-in”)

1. build empirical estimate \hat{P} for P
2. planning based on empirical \hat{P}

Model-based vs. model-free RL



Model-based approach (“plug-in”)

1. build empirical estimate \hat{P} for P
2. planning based on empirical \hat{P}

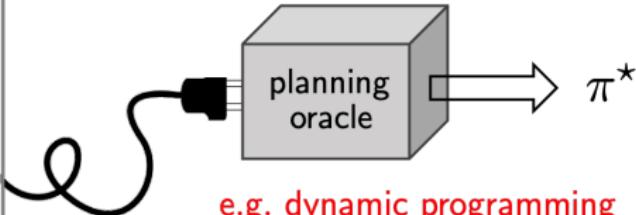
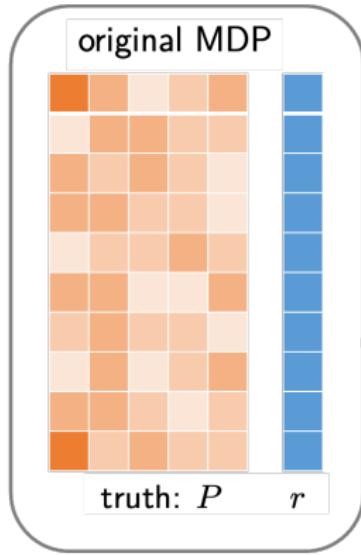
Model-free approach (e.g. Q-learning)

— learning w/o modeling & estimating environment explicitly

Vignette 1: Model-based RL (a “plug-in” approach)

“Breaking the sample size barrier in model-based reinforcement learning with a generative model,” G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, NeurIPS 2020

When the model is known ...

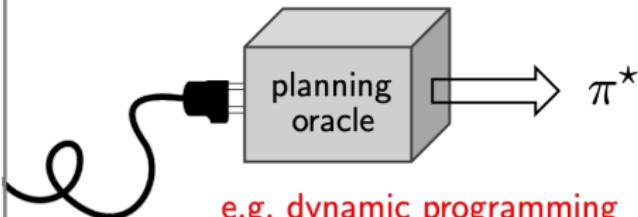
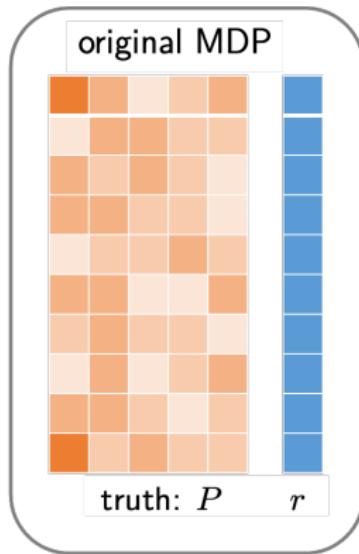


e.g. dynamic programming

1. Policy evaluation. Compute Q^{π_k}
2. Policy improvement. Update the policy: $\pi_{k+1} = \pi_{Q^{\pi_k}}$

Planning: computing the optimal policy π^* given the MDP specification

When the model is known ...



e.g. dynamic programming

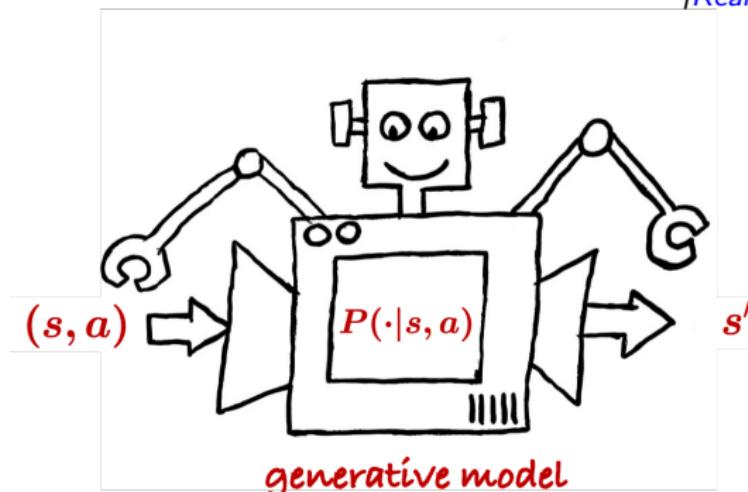
1. Policy evaluation. Compute Q^{π_k}
2. Policy improvement. Update the policy: $\pi_{k+1} = \pi_{Q^{\pi_k}}$

Planning: computing the optimal policy π^* given the MDP specification

In practice, do not know transition matrix P !

This work: sampling from a generative model

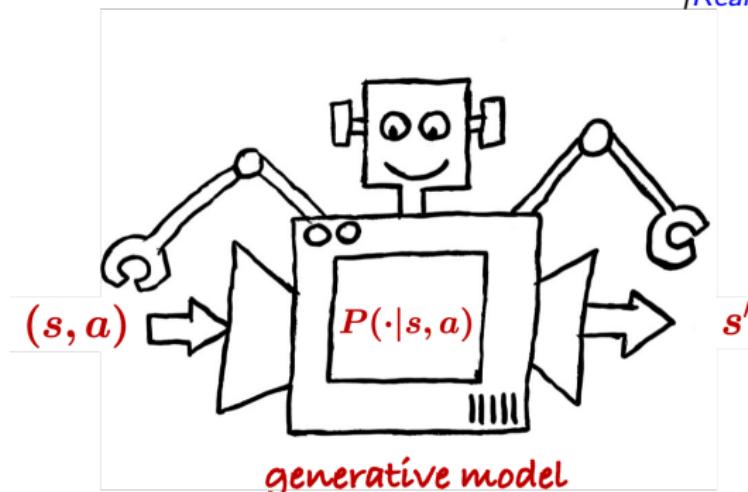
— [Kearns and Singh, 1999]



- **Sampling:** for each (s, a) , collect N samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

This work: sampling from a generative model

— [Kearns and Singh, 1999]



- **Sampling:** for each (s, a) , collect N samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$
- construct $\hat{\pi}$ based on samples (in total $|\mathcal{S}||\mathcal{A}| \times N$)

ℓ_∞ -sample complexity: how many samples are required to
learn an ε -optimal policy ?

$$\forall s: V^{\hat{\pi}}(s) \geq V^*(s) - \varepsilon$$

An incomplete list of prior art

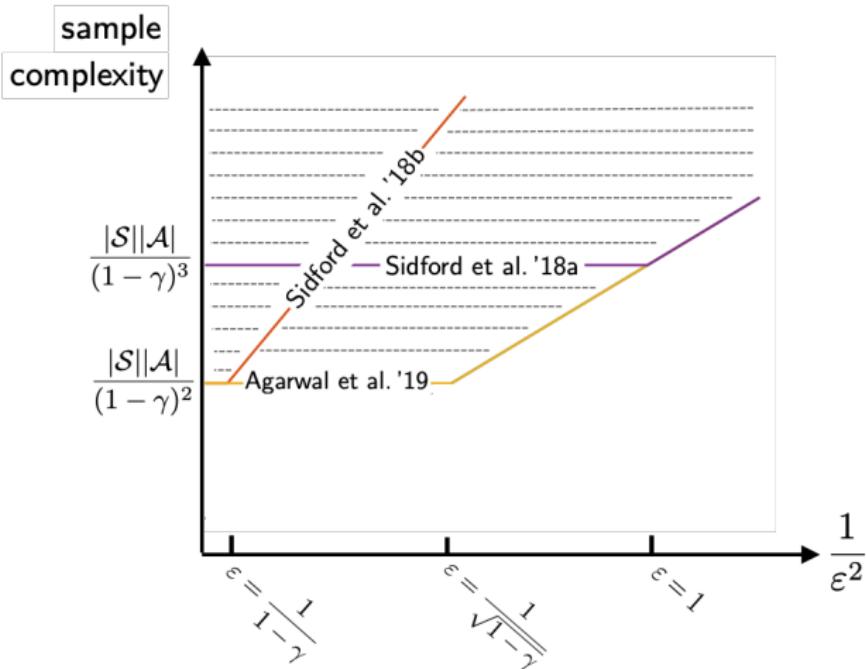
- [Kearns and Singh, 1999]
- [Kakade, 2003]
- [Kearns et al., 2002]
- [Azar et al., 2012]
- [Azar et al., 2013]
- [Sidford et al., 2018a]
- [Sidford et al., 2018b]
- [Wang, 2019]
- [Agarwal et al., 2019]
- [Wainwright, 2019a]
- [Wainwright, 2019b]
- [Pananjady and Wainwright, 2019]
- [Yang and Wang, 2019]
- [Khamaru et al., 2020]
- [Mou et al., 2020]
- ...

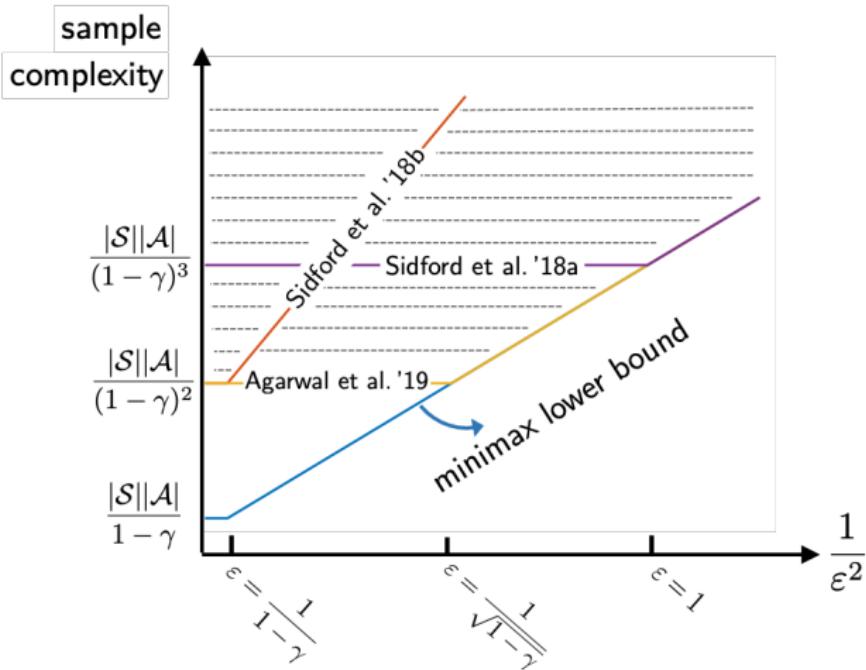
An even shorter list of prior art

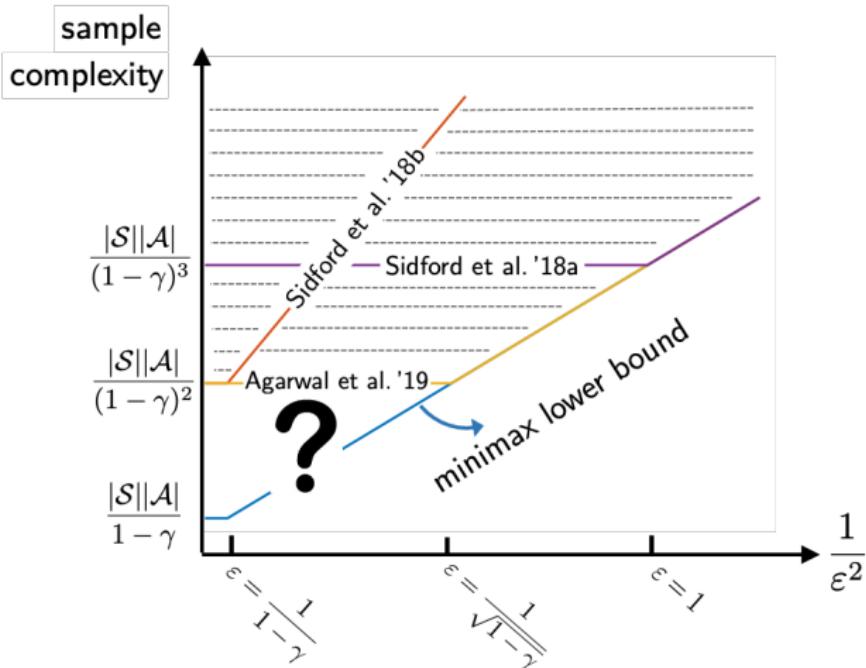
algorithm	sample size range	sample complexity	ε -range
Empirical QVI [Azar et al., 2013]	$[\frac{ \mathcal{S} ^2 \mathcal{A} }{(1-\gamma)^2}, \infty)$	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3 \varepsilon^2}$	$(0, \frac{1}{\sqrt{(1-\gamma) \mathcal{S} }}]$
Sublinear randomized VI [Sidford et al., 2018b]	$[\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^2}, \infty)$	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^4 \varepsilon^2}$	$(0, \frac{1}{1-\gamma}]$
Variance-reduced QVI [Sidford et al., 2018a]	$[\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3}, \infty)$	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3 \varepsilon^2}$	$(0, 1]$
Randomized primal-dual [Wang, 2019]	$[\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^2}, \infty)$	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^4 \varepsilon^2}$	$(0, \frac{1}{1-\gamma}]$
Empirical MDP + planning [Agarwal et al., 2019]	$[\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^2}, \infty)$	$\frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3 \varepsilon^2}$	$(0, \frac{1}{\sqrt{1-\gamma}}]$

important parameters:

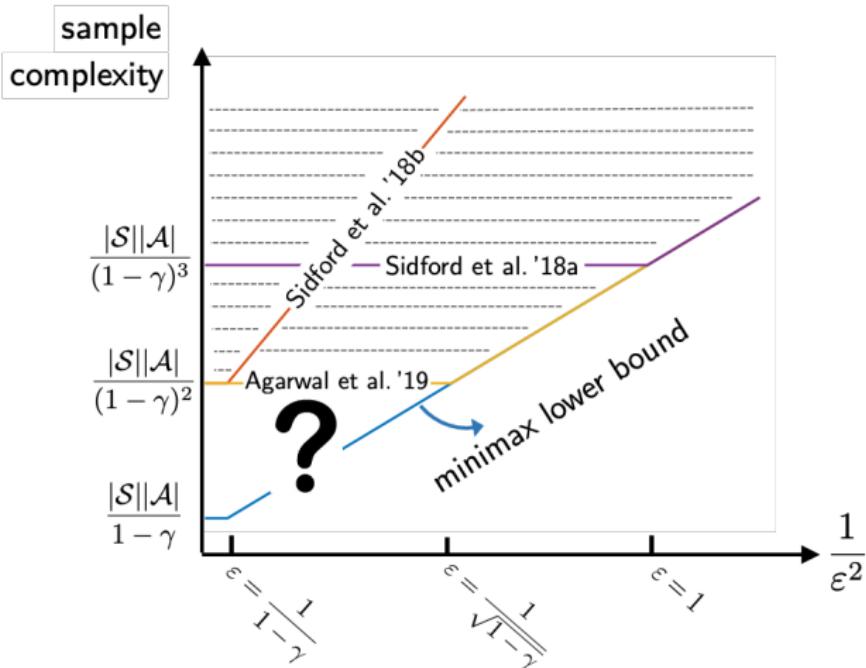
- $|\mathcal{S}|$: # states , $|\mathcal{A}|$: # actions
- $\frac{1}{1-\gamma}$: effective horizon
- $\varepsilon \in [0, \frac{1}{1-\gamma}]$: approximation error







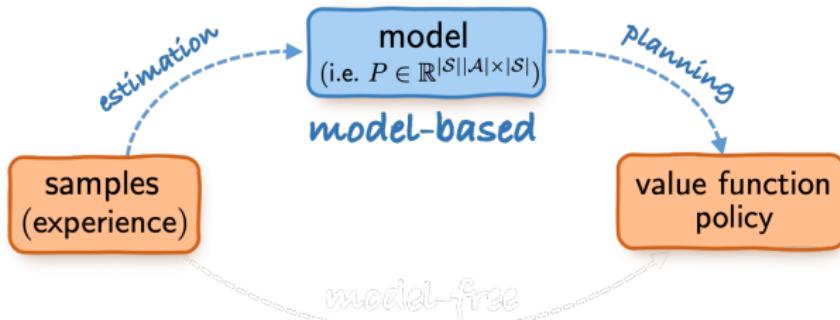
All prior theory requires sample size $\gtrsim \frac{|S||\mathcal{A}|}{(1-\gamma)^2}$



All prior theory requires sample size $\gtrsim \frac{|S||\mathcal{A}|}{(1-\gamma)^2}$

Question: is it possible to break this sample size barrier?

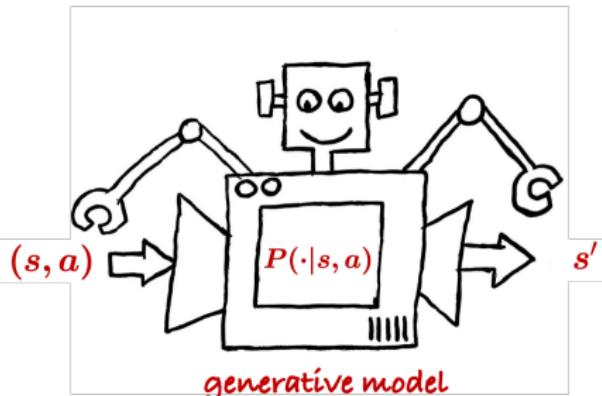
Our algorithm: model-based RL



Model-based approach (“plug-in”)

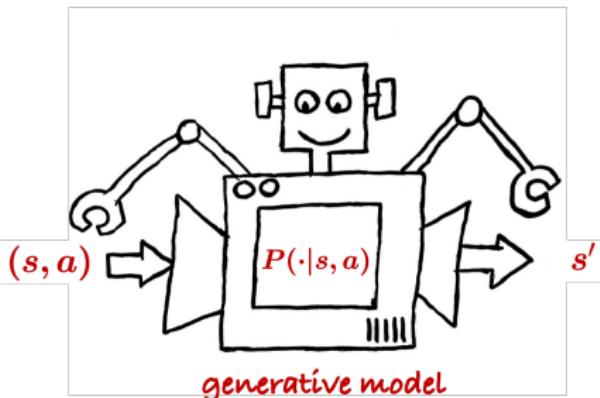
1. build an empirical estimate \hat{P} for P
2. planning based on empirical \hat{P}

Model estimation



Sampling: for each (s, a) , collect N ind. samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

Model estimation

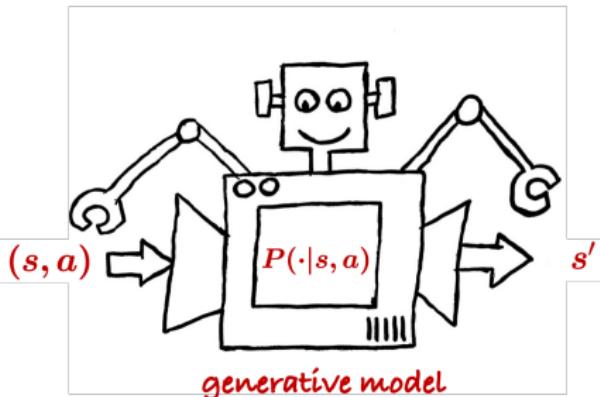


Sampling: for each (s, a) , collect N ind. samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

Empirical estimates:

$$\hat{P}(s'|s, a) = \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbb{1}\{s'_{(i)} = s'\}}_{\text{empirical frequency}}$$

Model estimation



Sampling: for each (s, a) , collect N ind. samples $\{(s, a, s'_{(i)})\}_{1 \leq i \leq N}$

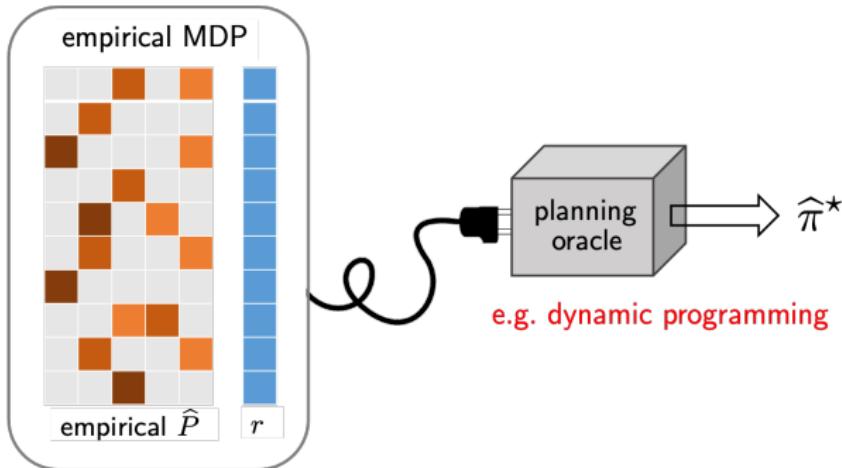
Empirical estimates:

$$\hat{P}(s'|s, a) = \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbb{1}\{s'_{(i)} = s'\}}_{\text{empirical frequency}}$$

- If sample size $\ll |\mathcal{S}|^2|\mathcal{A}|$, then we cannot recover P faithfully.

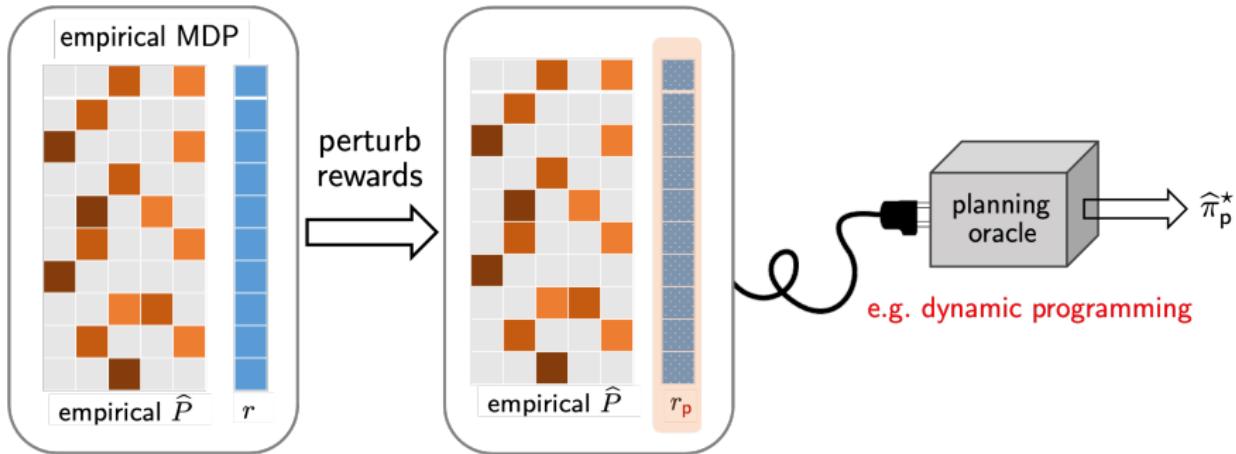
Model-based (plug-in) estimator

—[Azar et al., 2013, Agarwal et al., 2019, Pananjady and Wainwright, 2019]



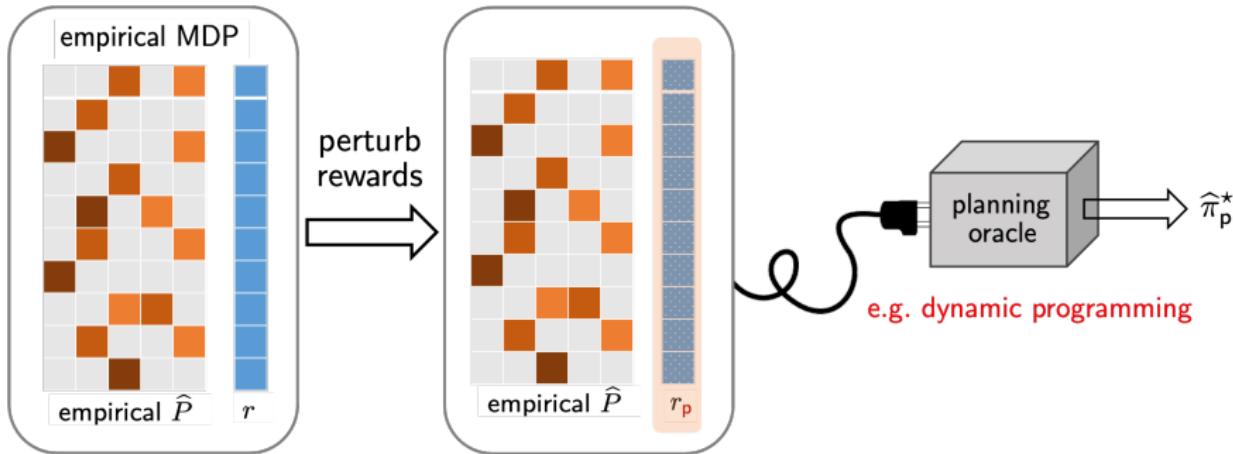
Find policy based on the **empirical** MDP (*empirical maximizer*)

Our method: plug-in estimator + perturbation



Find policy based on the **empirical** MDP with **slightly** perturbed rewards

Our method: plug-in estimator + perturbation



Find policy based on the **empirical** MDP with **slightly** perturbed rewards

Question: Can we trust our $\hat{\pi}$ when \hat{P} is not accurate?

Main result: ℓ_∞ -based sample complexity

Theorem (Li, Wei, Chi, Gu, Chen '20)

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the optimal policy $\widehat{\pi}_p^*$ of perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

Main result: ℓ_∞ -based sample complexity

Theorem (Li, Wei, Chi, Gu, Chen '20)

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the optimal policy $\widehat{\pi}_p^*$ of perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- minimax lower bound: $\widetilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$ [Azar et al., 2013]

Main result: ℓ_∞ -based sample complexity

Theorem (Li, Wei, Chi, Gu, Chen '20)

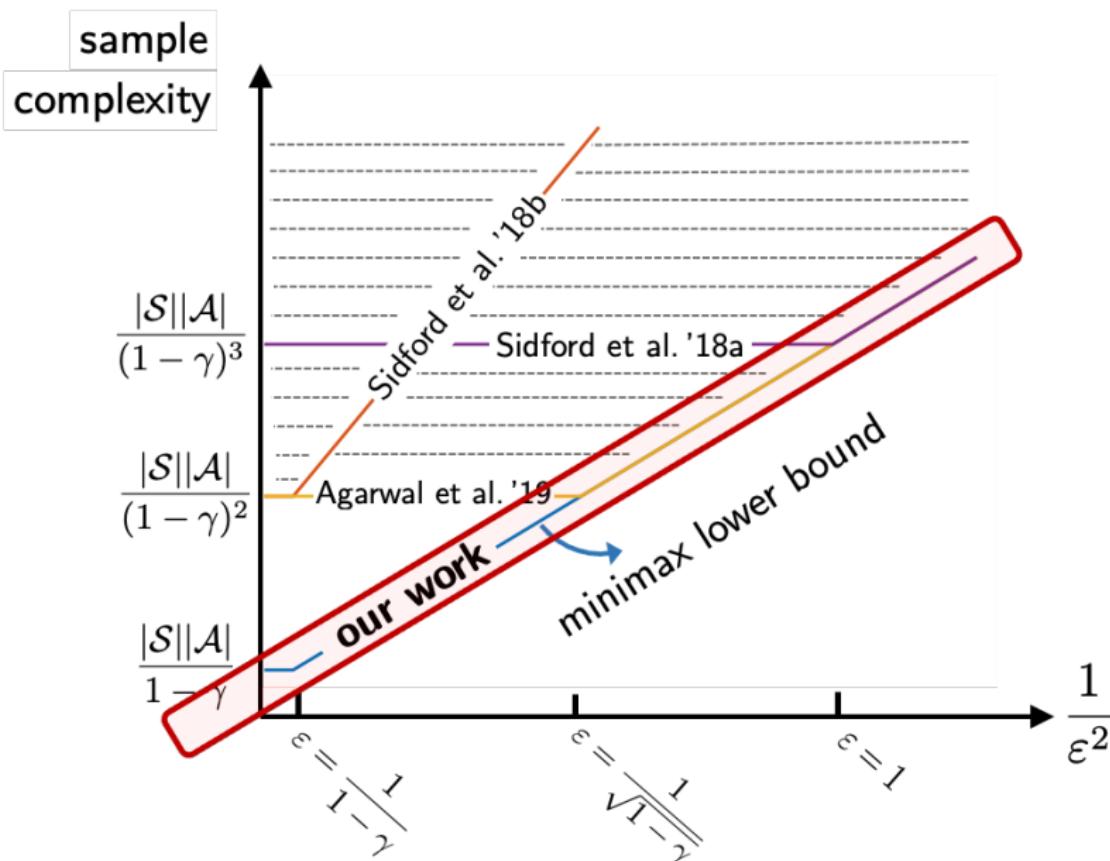
For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the optimal policy $\widehat{\pi}_p^*$ of perturbed empirical MDP achieves

$$\|V^{\widehat{\pi}_p^*} - V^*\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$$

- minimax lower bound: $\widetilde{\Omega}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3\varepsilon^2}\right)$ [Azar et al., 2013]
- $\varepsilon \in \left(0, \frac{1}{1-\gamma}\right] \rightarrow \text{sample size range } \left[\frac{|\mathcal{S}||\mathcal{A}|}{1-\gamma}, \infty\right)$



A glimpse of the key analysis ideas

Notation and Bellman equation

Bellman equation: $V^\pi = r_\pi + \gamma P_\pi V^\pi$

- V^π : value function under policy π
 - ▶ Bellman equation: $V^\pi = (I - \gamma P_\pi)^{-1} r_\pi$
- \hat{V}^π : empirical version value function under policy π
 - ▶ Bellman equation: $\hat{V}^\pi = (I - \gamma \hat{P}_\pi)^{-1} r_\pi$

Notation and Bellman equation

Bellman equation: $V^\pi = r_\pi + \gamma P_\pi V^\pi$

- V^π : value function under policy π
 - ▶ Bellman equation: $V^\pi = (I - \gamma P_\pi)^{-1} r_\pi$
- \hat{V}^π : empirical version value function under policy π
 - ▶ Bellman equation: $\hat{V}^\pi = (I - \gamma \hat{P}_\pi)^{-1} r_\pi$
- π^* : optimal policy for V^π
- $\hat{\pi}^*$: optimal policy for \hat{V}^π

Main steps

Elementary decomposition:

$$\begin{aligned} V^* - V^{\widehat{\pi}^*} &= (V^* - \widehat{V}^{\pi^*}) + (\widehat{V}^{\pi^*} - \widehat{V}^{\widehat{\pi}^*}) + (\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}) \\ &\leq (V^{\pi^*} - \widehat{V}^{\pi^*}) + \mathbf{0} + (\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}) \end{aligned}$$

Main steps

Elementary decomposition:

$$\begin{aligned} V^* - V^{\hat{\pi}^*} &= (V^* - \hat{V}^{\pi^*}) + (\hat{V}^{\pi^*} - \hat{V}^{\hat{\pi}^*}) + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \\ &\leq (V^{\pi^*} - \hat{V}^{\pi^*}) + \mathbf{0} + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \end{aligned}$$

- **Step 1:** control $V^\pi - \hat{V}^\pi$ for a fixed π (called “policy evaluation”)
(Bernstein inequality + a peeling argument)

Main steps

Elementary decomposition:

$$\begin{aligned} V^* - V^{\hat{\pi}^*} &= (V^* - \hat{V}^{\pi^*}) + (\hat{V}^{\pi^*} - \hat{V}^{\hat{\pi}^*}) + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \\ &\leq (V^{\pi^*} - \hat{V}^{\pi^*}) + \textcolor{red}{0} + (\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}) \end{aligned}$$

- **Step 1:** control $V^\pi - \hat{V}^\pi$ for a fixed π (called “policy evaluation”)
(Bernstein inequality + a peeling argument)
- **Step 2:** extend it to control $\hat{V}^{\hat{\pi}^*} - V^{\hat{\pi}^*}$ ($\hat{\pi}^*$ depends on samples)
(decouple statistical dependency)

Key idea 1: a peeling argument (for fixed policy)

First-order expansion

$$\hat{V}^\pi - V^\pi = \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)\hat{V}^\pi \quad [\text{Agarwal et al., 2019}]$$

Key idea 1: a peeling argument (for fixed policy)

First-order expansion

$$\hat{V}^\pi - V^\pi = \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)\hat{V}^\pi \quad [\text{Agarwal et al., 2019}]$$

Ours: higher-order expansion + Bernstein \longrightarrow tighter control

$$\begin{aligned}\hat{V}^\pi - V^\pi &= \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)V^\pi + \\ &\quad + \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)(\hat{V}^\pi - V^\pi)\end{aligned}$$

Key idea 1: a peeling argument (for fixed policy)

First-order expansion

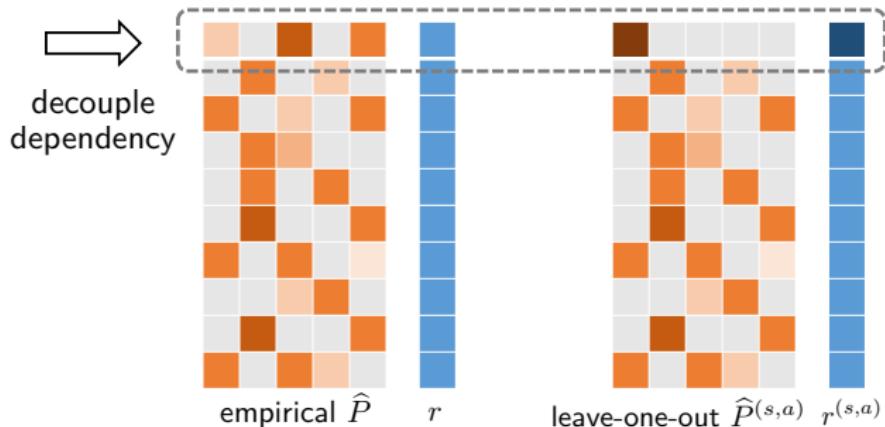
$$\hat{V}^\pi - V^\pi = \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)\hat{V}^\pi \quad [\text{Agarwal et al., 2019}]$$

Ours: higher-order expansion + Bernstein \longrightarrow tighter control

$$\begin{aligned}\hat{V}^\pi - V^\pi &= \gamma(I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi)V^\pi + \\ &\quad + \gamma^2 \left((I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi) \right)^2 V^\pi \\ &\quad + \gamma^3 \left((I - \gamma P_\pi)^{-1}(\hat{P}_\pi - P_\pi) \right)^3 V^\pi \\ &\quad + \dots\end{aligned}$$

Key idea 2: decouple dependency for $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$

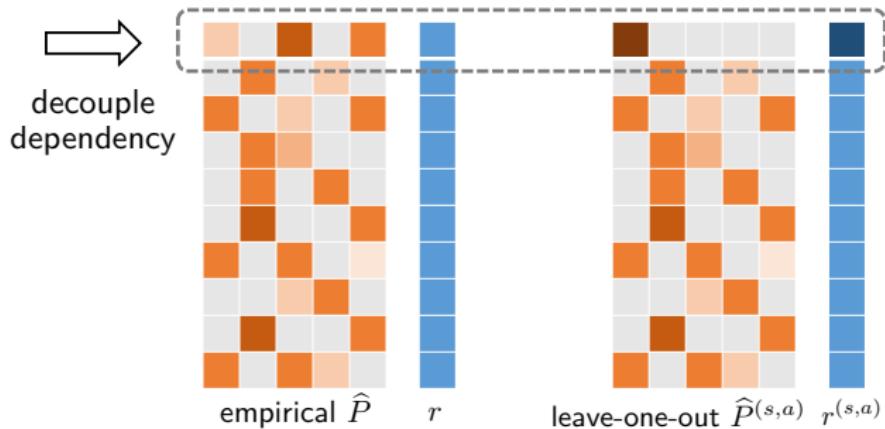
— inspired by [Agarwal et al., 2019] but quite different ...



- define $\widehat{\pi}_{(s,a)}^*$ $\xrightarrow{\text{empirical maximizer}} (\widehat{P}^{(s,a)}, r^{(s,a)})$

Key idea 2: decouple dependency for $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$

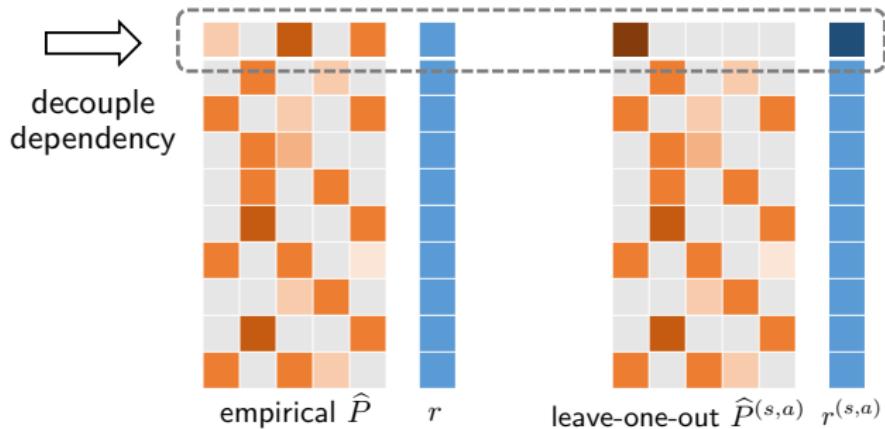
— inspired by [Agarwal et al., 2019] but quite different ...



- define $\widehat{\pi}_{(s,a)}^*$ $\xrightarrow{\text{empirical maximizer}} (\widehat{P}^{(s,a)}, r^{(s,a)})$
 - ▶ decouple dependency by dropping randomness in $\widehat{P}(\cdot | s, a)$
 - ▶ scalar $r^{(s,a)}$ ensures \widehat{Q}^* and \widehat{V}^* unchanged

Key idea 2: decouple dependency for $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$

— inspired by [Agarwal et al., 2019] but quite different ...



- define $\widehat{\pi}_{(s,a)}^*$ $\xrightarrow{\text{empirical maximizer}} (\widehat{P}^{(s,a)}, r^{(s,a)})$
- $\widehat{\pi}_{(s,a)}^* = \widehat{\pi}^*$ can be determined under separation condition

$$\forall s \in \mathcal{S}, \quad \widehat{Q}^*(s, \widehat{\pi}^*(s)) - \max_{a: a \neq \widehat{\pi}^*(s)} \widehat{Q}^*(s, a) > 0$$

Key idea 3: tie-breaking via reward perturbation

- How to ensure separation btw the optimal policy and others?

$$\forall s \in \mathcal{S}, \quad \widehat{Q}^*(s, \widehat{\pi}^*(s)) - \max_{a: a \neq \widehat{\pi}^*(s)} \widehat{Q}^*(s, a) > 0$$

Key idea 3: tie-breaking via reward perturbation

- How to ensure separation btw the optimal policy and others?

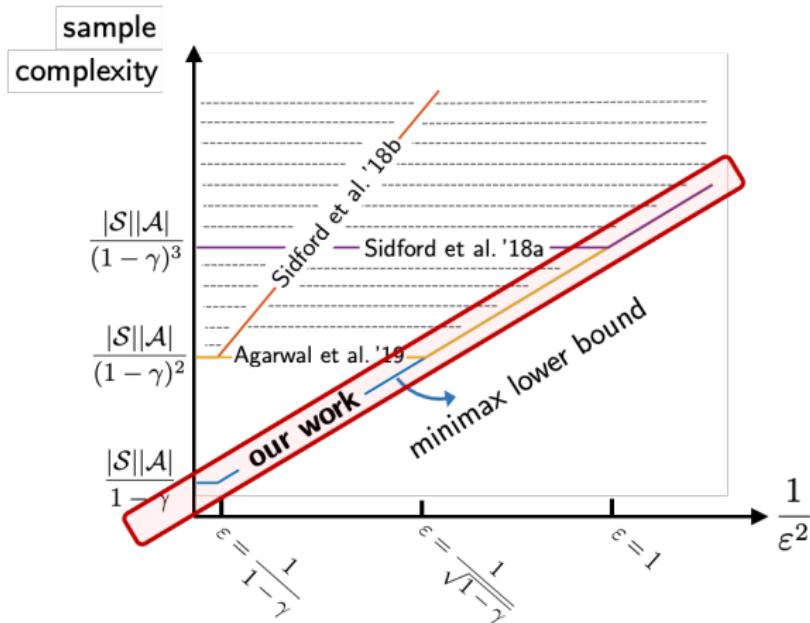
$$\forall s \in \mathcal{S}, \quad \widehat{Q}^*(s, \widehat{\pi}^*(s)) - \max_{a: a \neq \widehat{\pi}^*(s)} \widehat{Q}^*(s, a) > 0$$

- Solution:** slightly perturb rewards $r \implies \widehat{\pi}_p^*$

- ensures $\widehat{\pi}_p^*$ can be differentiated from others
- $V^{\widehat{\pi}_p^*} \approx V^{\widehat{\pi}^*}$



Summary of this part

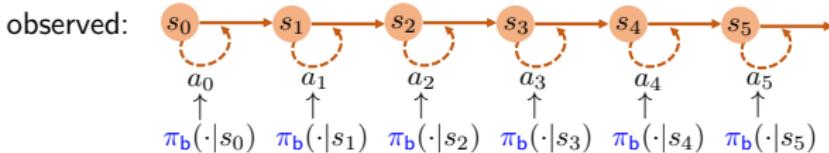


Model-based RL is minimax optimal & does not suffer from a sample size barrier!

Vignette 2: Model-free approach

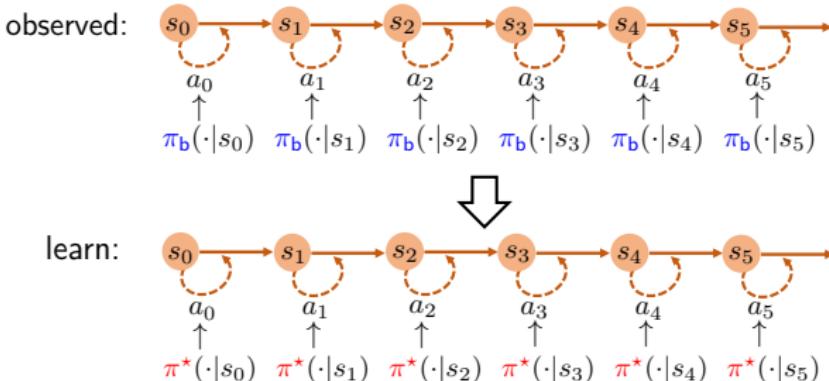
"Sample Complexity of Asynchronous Q-Learning: Sharper Analysis and Variance Reduction," G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, NeurIPS 2020

Markovian samples and behavior policy



Observed: $\underbrace{\{s_t, a_t, r_t\}_{t \geq 0}}_{\text{Markovian trajectory}}$ induced by behavior policy π_b

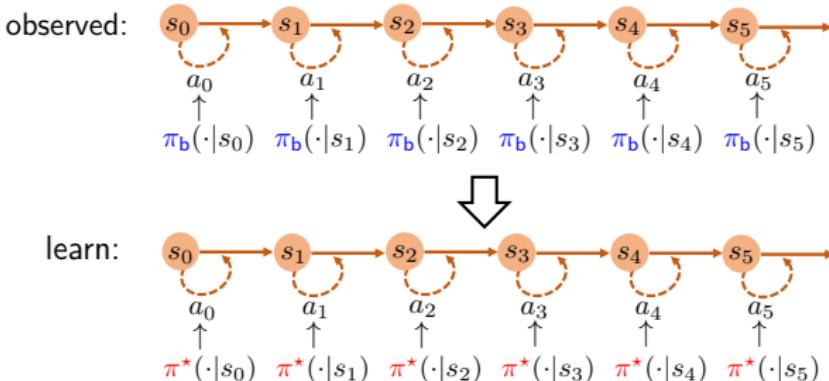
Markovian samples and behavior policy



Observed: $\underbrace{\{s_t, a_t, r_t\}_{t \geq 0}}_{\text{Markovian trajectory}}$ induced by **behavior policy** π_b

Goal: learn optimal value V^* and Q^* based on sample trajectory

Markovian samples and behavior policy



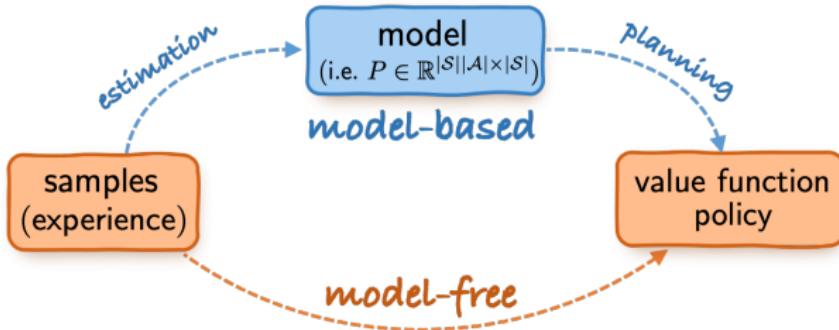
Key quantities of sample trajectory

- minimum state-action occupancy probability

$$\mu_{\min} := \min \underbrace{\mu_{\pi_b}(s, a)}_{\text{stationary distribution}}$$

- mixing time: t_{mix}

Model-based vs. model-free RL



Model-free approach (e.g. Q-learning)

— learning w/o modeling & estimating environment explicitly

Q-learning: a classical model-free algorithm



Chris Watkins



Peter Dayan

Stochastic approximation for solving **Bellman equation** $Q = \mathcal{T}(Q)$

Robbins & Monro '51

Aside: Bellman optimality principle

Bellman operator

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right]$$

- one-step look-ahead

Aside: Bellman optimality principle

Bellman operator

$$\mathcal{T}(Q)(s, a) := \underbrace{r(s, a)}_{\text{immediate reward}} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\underbrace{\max_{a' \in \mathcal{A}} Q(s', a')}_{\text{next state's value}} \right]$$

- one-step look-ahead

Bellman equation: Q^* is *unique* solution to

$$\mathcal{T}(Q^*) = Q^*$$



Richard Bellman

Q-learning: a classical model-free algorithm



Chris Watkins



Peter Dayan

Stochastic approximation for solving Bellman equation $Q = \mathcal{T}(Q)$

$$\underbrace{Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \eta_t(\mathcal{T}_t(Q_t)(s_t, a_t) - Q_t(s_t, a_t))}_{\text{only update } (s_t, a_t)\text{-th entry}}, \quad t \geq 0$$

Q-learning: a classical model-free algorithm



Chris Watkins



Peter Dayan

Stochastic approximation for solving Bellman equation $Q = \mathcal{T}(Q)$

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \eta_t (\underbrace{\mathcal{T}_t(Q_t)(s_t, a_t) - Q_t(s_t, a_t)}_{\text{only update } (s_t, a_t)\text{-th entry}}), \quad t \geq 0$$

$$\mathcal{T}_t(Q)(s_t, a_t) := r(s_t, a_t) + \gamma \max_{a'} Q(s_{t+1}, a')$$

$$\mathcal{T}(Q)(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a'} Q(s', a') \right]$$

Q-learning: a classical model-free algorithm



Chris Watkins



Peter Dayan

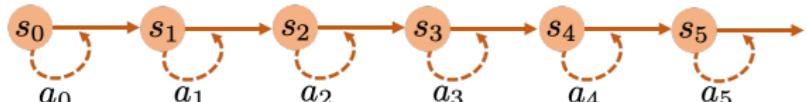
Stochastic approximation for solving Bellman equation $Q = \mathcal{T}(Q)$

$$Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \eta_t (\mathcal{T}_t(Q_t)(s_t, a_t) - Q_t(s_t, a_t)), \quad t \geq 0$$

only update (s_t, a_t) -th entry

— **asynchronous:** only a single entry is updated each iteration
(resembles Markov-chain *coordinate descent*)

observed:



What is sample complexity of (async) Q-learning?

A highly incomplete list of prior work

- [Watkins and Dayan, 1992]
- [Tsitsiklis, 1994]
- [Jaakkola et al., 1994]
- [Szepesvári, 1998]
- [Kearns and Singh, 1999]
- [Borkar and Meyn, 2000]
- [Even-Dar and Mansour, 2003]
- [Beck and Srikant, 2012]
- [Jin et al., 2018]
- [Shah and Xie, 2018]
- [Wainwright, 2019a]
- [Chen et al., 2019]
- [Yang and Wang, 2019]
- [Du et al., 2020]
- [Chen et al., 2020]
- [Qu and Wierman, 2020]
- [Devraj and Meyn, 2020]
- ...

Prior art: async Q-learning

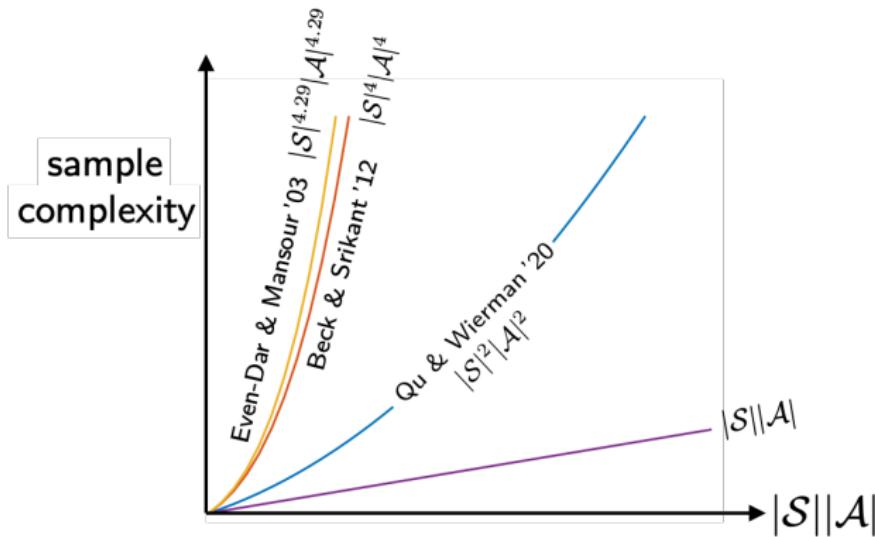
Question: how many samples are needed to ensure $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$?

paper	sample complexity	learning rate
[Even-Dar and Mansour, 2003]	$\frac{(t_{\text{cover}})^{\frac{1}{1-\gamma}}}{(1-\gamma)^4 \varepsilon^2}$	linear: $\frac{1}{t}$
[Even-Dar and Mansour, 2003]	$\left(\frac{t_{\text{cover}}^{1+3\omega}}{(1-\gamma)^4 \varepsilon^2}\right)^{\frac{1}{\omega}} + \left(\frac{t_{\text{cover}}}{1-\gamma}\right)^{\frac{1}{1-\omega}}$	poly: $\frac{1}{t^\omega}$, $\omega \in (\frac{1}{2}, 1)$
[Beck and Srikant, 2012]	$\frac{t_{\text{cover}}^3 \mathcal{S} \mathcal{A} }{(1-\gamma)^5 \varepsilon^2}$	constant
[Qu and Wierman, 2020]	$\frac{t_{\text{mix}}}{\mu_{\min}^2 (1-\gamma)^5 \varepsilon^2}$	rescaled linear

— cover time: $t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$

Prior art: async Q-learning

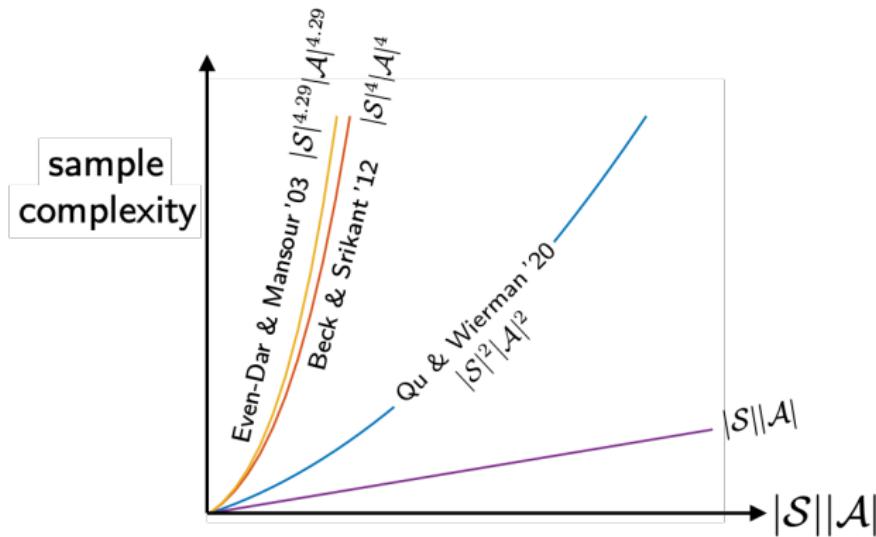
Question: how many samples are needed to ensure $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$?



if we take $\mu_{\min} \asymp \frac{1}{|S||\mathcal{A}|}$, $t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$

Prior art: async Q-learning

Question: how many samples are needed to ensure $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$?

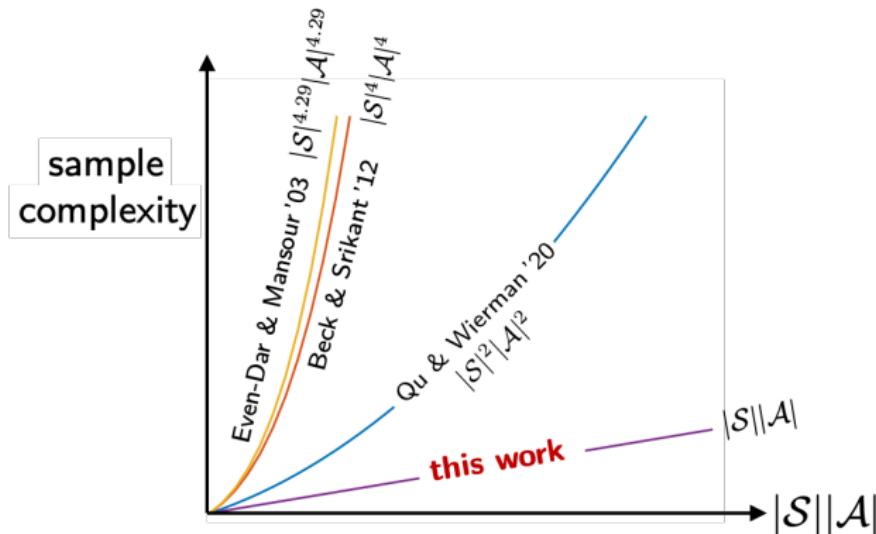


$$\text{if we take } \mu_{\min} \asymp \frac{1}{|S||A|}, t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$$

All prior results require sample size of at least $t_{\text{mix}}|S|^2|\mathcal{A}|^2$!

Prior art: async Q-learning

Question: how many samples are needed to ensure $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$?



if we take $\mu_{\min} \asymp \frac{1}{|S||\mathcal{A}|}$, $t_{\text{cover}} \asymp \frac{t_{\text{mix}}}{\mu_{\min}}$

All prior results require sample size of at least $t_{\text{mix}}|S|^2|\mathcal{A}|^2$!

Main result: ℓ_∞ -based sample complexity

Theorem (Li, Wei, Chi, Gu, Chen '20)

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, sample complexity of async Q-learning to yield $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ is at most (up to some log factor)

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

Main result: ℓ_∞ -based sample complexity

Theorem (Li, Wei, Chi, Gu, Chen '20)

For any $0 < \varepsilon \leq \frac{1}{1-\gamma}$, sample complexity of async Q-learning to yield $\|\hat{Q} - Q^*\|_\infty \leq \varepsilon$ is at most (up to some log factor)

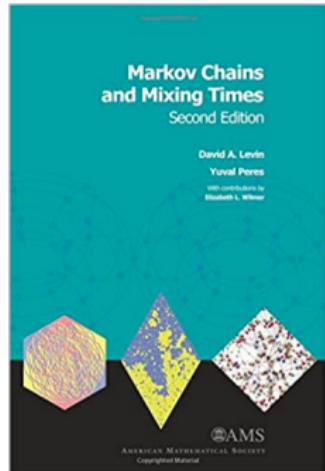
$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

— prior art: $\frac{t_{\text{mix}}}{\mu_{\min}^2(1-\gamma)^5 \varepsilon^2}$ ([Qu and Wierman, 2020])

- Improves upon prior art by **at least** $|\mathcal{S}||\mathcal{A}|!$

Effect of mixing time on sample complexity

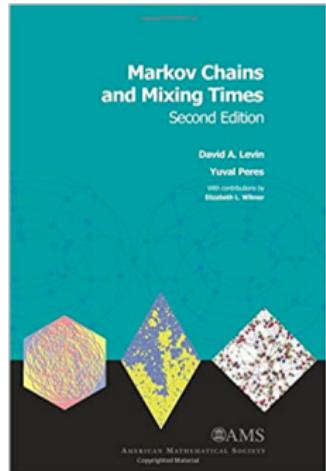
$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$



- reflects cost taken to reach steady state
- one-time expense (almost independent of ε)
 - it becomes amortized as algorithm runs

Effect of mixing time on sample complexity

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$



- reflects cost taken to reach steady state
- one-time expense (almost independent of ε)
 - it becomes amortized as algorithm runs

— prior art: $\frac{t_{\text{mix}}}{\mu_{\min}^2(1-\gamma)^5 \varepsilon^2}$ ([Qu and Wierman, 2020])

Dependence on effective horizon

minimax lower bound
(Azar et al. '13)

$$\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2}$$

asyn Q-learning
(ignoring dependency on t_{mix})

$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2}$$

Dependence on effective horizon

minimax lower bound
(Azar et al. '13)

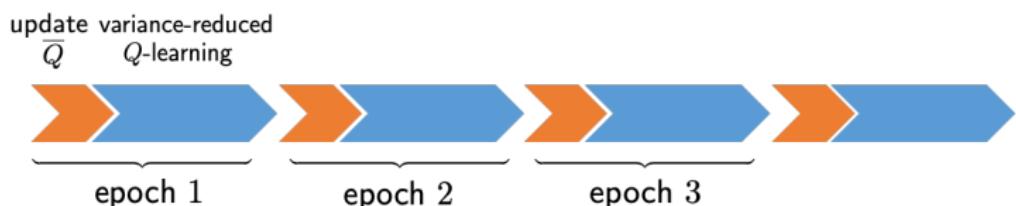
$$\frac{1}{\mu_{\min}(1-\gamma)^3 \varepsilon^2}$$

asyn Q-learning
(ignoring dependency on t_{mix})

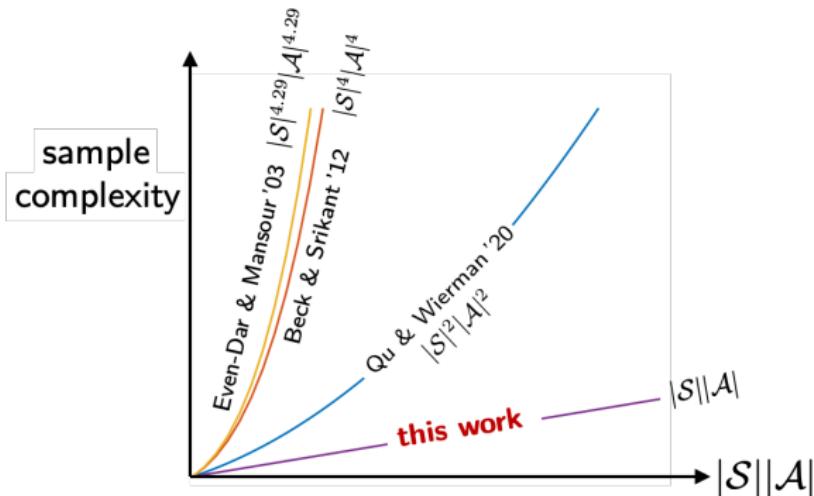
$$\frac{1}{\mu_{\min}(1-\gamma)^5 \varepsilon^2}$$

The dependency on $\frac{1}{1-\gamma}$ can be tightened by *variance reduction*.

— inspired by [[Johnson and Zhang, 2013](#)], [[Wainwright, 2019b](#)]



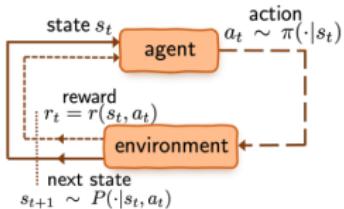
Summary of this part



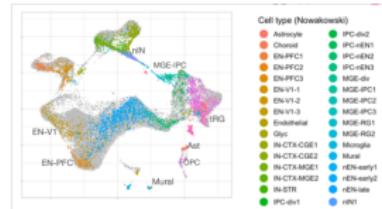
Sharper sample complexity for asyn Q-learning
in terms of $|S||\mathcal{A}|$ and t_{mix} !

Concluding remark

Integration of modern statistics, optimization and RL is beneficial.



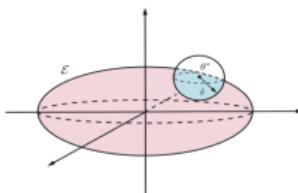
Reinforcement Learning



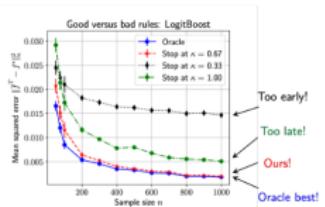
Transfer Learning for scRNASEq

$$y = \underbrace{\begin{matrix} p \\ X \end{matrix}}_n + z$$

Lasso Asymptotics



Shape Constrained Inference



Early Stopping for Boosting

Thanks for your attention!

Other details

Improved theory for policy evaluation

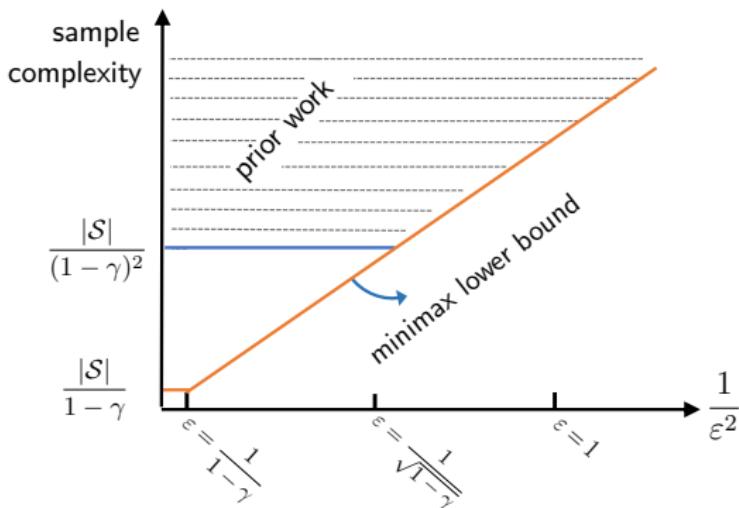
Model-based policy evaluation:

- given a fixed policy π , estimate V^π via the plug-in estimate \hat{V}^π

Improved theory for policy evaluation

Model-based policy evaluation:

- given a fixed policy π , estimate V^π via the plug-in estimate \hat{V}^π



- A sample size barrier $\frac{|\mathcal{S}|}{(1-\gamma)^2}$ already appeared in prior work
(Agarwal et al. '19, Pananjady & Wainwright '19, Khamaru et al. '20)

Improved theory for policy evaluation

Model-based policy evaluation:

- given a fixed policy π , estimate V^π via the plug-in estimate \hat{V}^π

Theorem (Li, Wei, Chi, Gu, Chen'20)

Fix any policy π . For $0 < \varepsilon \leq \frac{1}{1-\gamma}$, the plug-in estimator \hat{V}^π obeys

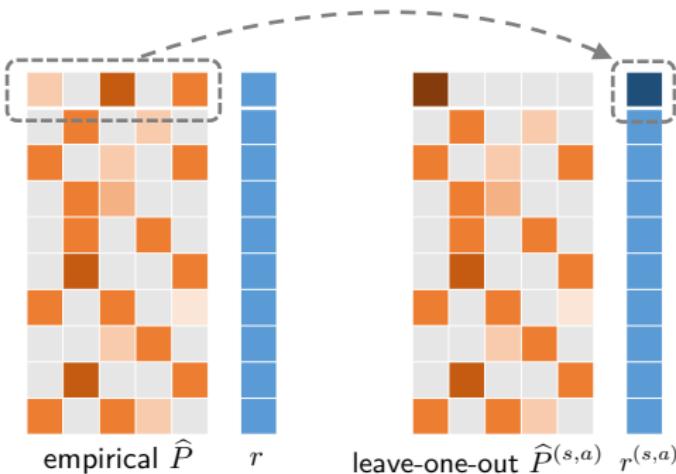
$$\|\hat{V}^\pi - V^\pi\|_\infty \leq \varepsilon$$

with sample complexity at most

$$\tilde{O}\left(\frac{|\mathcal{S}|}{(1-\gamma)^3 \varepsilon^2}\right)$$

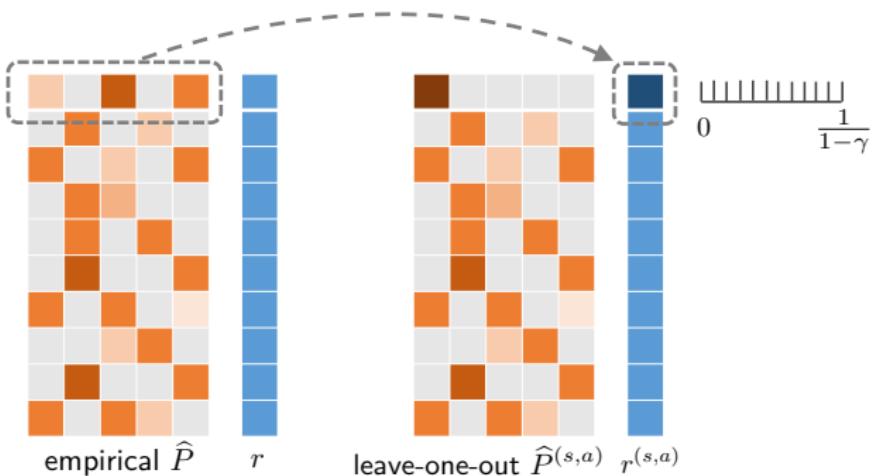
- Minimax optimal for all ε (Azar et al. '13, Pananjady & Wainwright '19)

Key idea 2: leave-one-out analysis for $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$



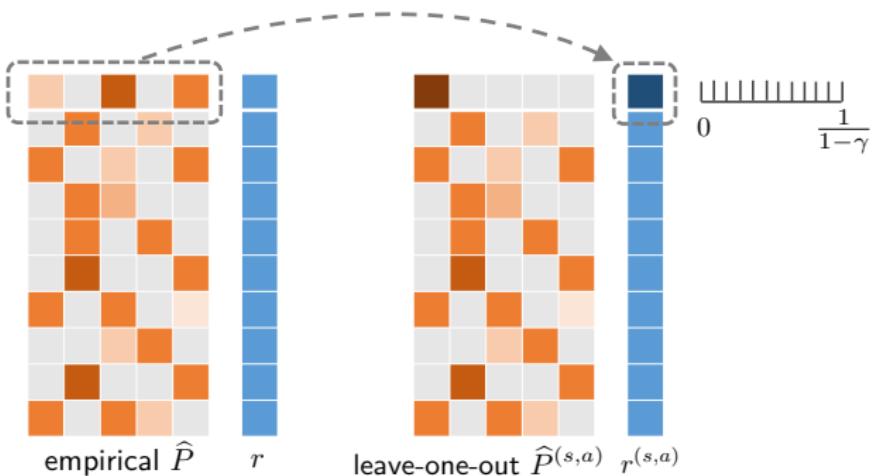
1. embed all randomness from $\widehat{P}(\cdot \mid s, a)$ into a single scalar (i.e. $r^{(s,a)}$)

Key idea 2: leave-one-out analysis for $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$



1. embed all randomness from $\widehat{P}(\cdot | s, a)$ into a single scalar (i.e. $r^{(s,a)}$)
2. build an ϵ -net for this scalar

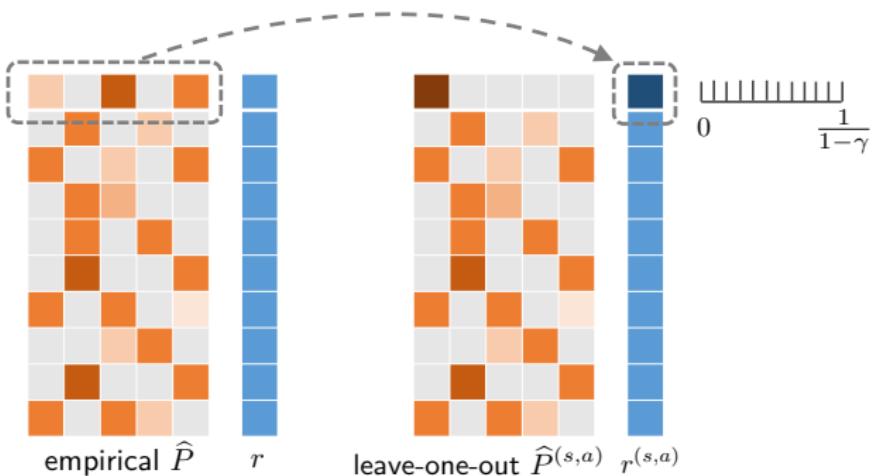
Key idea 2: leave-one-out analysis for $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$



1. embed all randomness from $\widehat{P}(\cdot | s, a)$ into a single scalar (i.e. $r^{(s,a)}$)
2. build an ϵ -net for this scalar
3. $\widehat{\pi}_{(s,a)}^* = \widehat{\pi}^*$ can be determined under separation condition

$$\forall s \in \mathcal{S}, \quad \widehat{Q}^*(s, \widehat{\pi}^*(s)) - \max_{a: a \neq \widehat{\pi}^*(s)} \widehat{Q}^*(s, a) > 0$$

Key idea 2: leave-one-out analysis for $\widehat{V}^{\widehat{\pi}^*} - V^{\widehat{\pi}^*}$



Compared to [Agarwal et al., 2019]

- [Agarwal et al., 2019]: dependency btw value \widehat{V} & samples
- **Ours:** dependency btw policy $\widehat{\pi}$ & samples

Key decomposition for asyn Q-learning

Error decomposition

$$\Delta_t = (\mathbf{I} - \Lambda_t) \Delta_{t-1} + \gamma \Lambda_t (\mathbf{P}_t - \mathbf{P}) \mathbf{V}^* + \gamma \Lambda_t \mathbf{P}_t (\mathbf{V}_{t-1} - \mathbf{V}^*)$$

Applying this relation recursively gives

$$\begin{aligned}\Delta_t &= \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i (\mathbf{P}_i - \mathbf{P}) \mathbf{V}^* \\ &\quad + \gamma \sum_{i=1}^t \prod_{j=i+1}^t (\mathbf{I} - \Lambda_j) \Lambda_i \mathbf{P}_i (\mathbf{V}_{i-1} - \mathbf{V}^*) + \prod_{j=1}^t (\mathbf{I} - \Lambda_j) \Delta_0\end{aligned}$$

Variance-reduced Q-learning

- Update for reference Bellman operator:

$$\tilde{\mathcal{T}}(\overline{Q})(s, a) = r(s, a) + \frac{\gamma \sum_{i=0}^{N-1} \mathbb{1}\{(s_i, a_i) = (s, a)\} \max_{a'} \overline{Q}(s_{i+1}, a')}{\sum_{i=0}^{N-1} \mathbb{1}\{(s_i, a_i) = (s, a)\}}$$

- Parameter choices:

$$\eta_t = \frac{c_0}{\log\left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta}\right)} \min\left\{\frac{(1-\gamma)^2}{\gamma^2}, \frac{1}{t_{\text{mix}}}\right\}$$

$$N = \frac{1}{\mu_{\min}} \left(\frac{1}{(1-\gamma)^3 \min\{1, \varepsilon^2\}} + t_{\text{mix}} \right) \log\left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta}\right)$$

$$t_{\text{epoch}} = \frac{1}{\mu_{\min}} \left(\frac{1}{(1-\gamma)^3} + \frac{t_{\text{mix}}}{1-\gamma} \right) \log\left(\frac{1}{(1-\gamma)^2 \varepsilon}\right) \log\left(\frac{|\mathcal{S}||\mathcal{A}|t_{\text{epoch}}}{\delta}\right)$$

Learning rates

$$\text{constant stepsize } \eta_t \equiv \min \left\{ \frac{(1-\gamma)^4 \varepsilon^2}{\gamma^2}, \frac{1}{t_{\text{mix}}} \right\}$$

- [Qu and Wierman, 2020]: rescaled linear $\eta_t = \frac{\frac{1}{\mu_{\min}(1-\gamma)}}{t + \max\{\frac{1}{\mu_{\min}(1-\gamma)}, t_{\text{mix}}\}}$
- [Beck and Srikant, 2012] constant $\eta_t \equiv \underbrace{\frac{(1-\gamma)^4 \varepsilon^2}{|\mathcal{S}||\mathcal{A}|t_{\text{cover}}^2}}_{\text{too conservative}}$
- [Even-Dar and Mansour, 2003]: polynomial $\eta_t = t^{-\omega}$ ($\omega \in (\frac{1}{2}, 1]$)

Adaptive learning rates

$$\eta_t = \min \left\{ 1, c \exp \left(\left| \log \frac{\log t}{\hat{\mu}_{\min,t} (1 - \gamma) \gamma^2 t} \right| \right) \right\}$$

$$\hat{\mu}_{\min,t} = \begin{cases} \frac{1}{|\mathcal{S}||\mathcal{A}|}, & \min_{s,a} K_t(s,a) = 0; \\ \hat{\mu}_{\min,t-1}, & \frac{1}{2} < \frac{\min_{s,a} K_t(s,a)/t}{\hat{\mu}_{\min,t-1}} < 2; \\ \min_{s,a} K_t(s,a)/t, & \text{otherwise.} \end{cases}$$

One strategy: variance reduction

— inspired by [Johnson and Zhang, 2013], [Wainwright, 2019b]

Variance-reduced Q-learning updates

$$Q_t(s_t, a_t) = (1 - \eta)Q_{t-1}(s_t, a_t) + \eta \left(\mathcal{T}_t(Q_{t-1}) \underbrace{- \mathcal{T}_t(\bar{Q}) + \tilde{\mathcal{T}}(\bar{Q})}_{\text{use } \bar{Q} \text{ to help reduce variability}} \right)(s_t, a_t)$$

- \bar{Q} : some reference Q-estimate
- $\tilde{\mathcal{T}}$: empirical Bellman operator (using a batch of samples)

Variance-reduced Q-learning

— inspired by [[Johnson and Zhang, 2013](#)], [[Wainwright, 2019b](#)]

update variance-reduced
 \bar{Q} Q -learning



for each epoch

1. update \bar{Q} and $\tilde{T}(\bar{Q})$
2. run variance-reduced Q-learning updates

Main result: ℓ_∞ -based sample complexity

Theorem (Li, Wei, Chi, Gu, Chen '20)

For any $0 < \varepsilon \leq 1$, sample complexity for (async) variance-reduced Q-learning to yield $\|\widehat{Q} - Q^*\|_\infty \leq \varepsilon$ is at most on the order of

$$\frac{1}{\mu_{\min}(1-\gamma)^3\varepsilon^2} + \frac{t_{\text{mix}}}{\mu_{\min}(1-\gamma)}$$

- more aggressive learning rates: $\eta_t \equiv \min \left\{ \frac{(1-\gamma)^4(1-\gamma)^2}{\gamma^2}, \frac{1}{t_{\text{mix}}} \right\}$
- minimax-optimal for $0 < \varepsilon \leq 1$